

DOCUMENT PROCESSOR AND STORAGE MEDIUM

Publication number: JP11259498

Publication date: 1999-09-24

Inventor: HASHIMOTO MINAKO; KASHINO WAKAKO;
OCHITANI AKIRA; NISHINO FUMITO

Applicant: FUJITSU LTD

Classification:

- International: G06F17/21; G06F17/30; G06F17/21; G06F17/30;
(IPC1-7): G06F17/30; G06F17/21

- European: G06F17/30T1E

Application number: JP19980058384 19980310

Priority number(s): JP19980058384 19980310

Also published as:

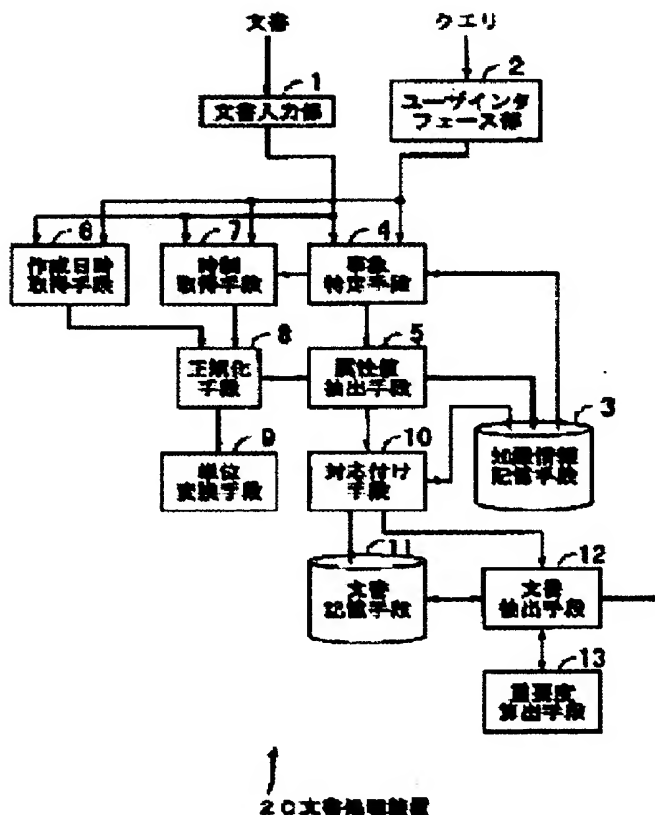
WO9946698 (A1)
US6523025 (B1)
GB2350712 (A)
CN1286776 (A)
CN1109994C (C)

Report a data error here

Abstract of JP11259498

PROBLEM TO BE SOLVED: To enhance retrieving or clipping accuracy of a document.

SOLUTION: The document to be processed is inputted from a document input part 1. A type of an event described on the inputted document is specified by referring to knowledge information stored in a knowledge information storage means 3 by an event specifying means 4. An attribute value of an attribute regarding the specified event is extracted from the document by an attribute value extracting means 5. A processing to make the attribute value extracted by the attribute value extracting means 5 correspond to substance in the real world is performed by a corresponding means 10. Information (normalized information) generated by the corresponding means 10 and information to specify the document or a place in which the document is stored are stored by correlating them by a document storage means 11. A query inputted from a user interface part 2 is compared with the normalized information and the information to specify the applicable document or the place in which the document is stored is outputted from the document storage means 11 by a document extracting means 12.



Data supplied from the esp@cenet database - Worldwide

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開平11-259498

(43)公開日 平成11年(1999) 9月24日

(51)Int.Cl.⁸

識別記号

F I

G 0 6 F 17/30
17/21

G 0 6 F 15/40 3 7 0 A
15/20 5 9 0 E
15/403 3 1 0 A
3 5 0 C

審査請求 未請求 請求項の数10 O L (全 20 頁)

(21)出願番号 特願平10-58384

(22)出願日 平成10年(1998) 3月10日

(71)出願人 000005223

富士通株式会社

神奈川県川崎市中原区上小田中4丁目1番
1号

(72)発明者 橋本 三奈子

神奈川県川崎市中原区上小田中4丁目1番
1号 富士通株式会社内

(72)発明者 柏野 和佳子

神奈川県川崎市中原区上小田中4丁目1番
1号 富士通株式会社内

(74)代理人 弁理士 服部 毅巖

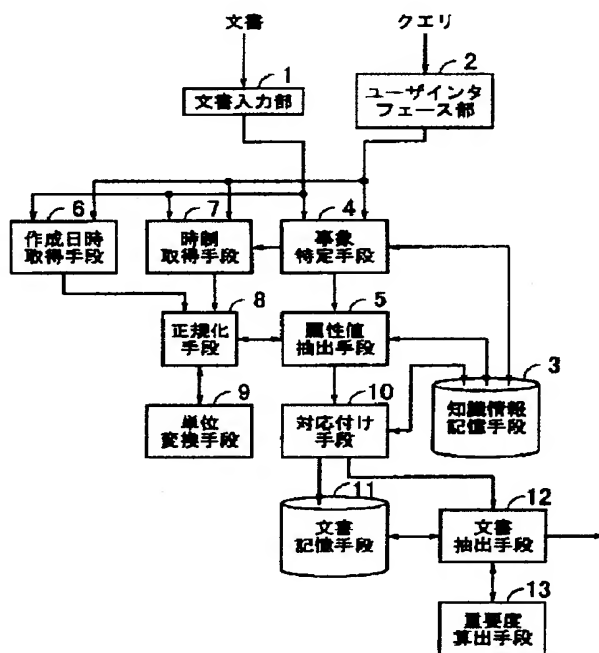
最終頁に続く

(54)【発明の名称】 文書処理装置および記録媒体

(57)【要約】

【課題】 文書の検索またはクリッピング精度を向上させる。

【解決手段】 文書入力部1からは処理の対象となる文書が入力される。事象特定手段4は、知識情報記憶手段3に記憶されている知識情報を参照して、入力された文書に記述されている事象の種類を特定する。属性値抽出手段5は、特定された事象に係わる属性の属性値を文書から抽出する。対応付け手段10は、属性値抽出手段5によって抽出された属性値を、実世界の実体に対して対応付ける処理を行う。文書記憶手段11は、対応付け手段10によって生成された情報(正規化情報)と、文書または文書の格納場所を特定する情報とを関連付けて記憶する。文書抽出手段12は、ユーザインタフェース部2から入力されたクエリと、正規化情報とを比較し、該当する文書または文書の格納場所を特定する情報を文書記憶手段11から出力する。



↑
20 文書処理装置

【特許請求の範囲】

【請求項 1】 入力された文書に対して所定の処理を施して記憶するとともに、与えられたクエリに対応する文書を、記憶されている文書の中から検索またはクリッピングする処理を行う文書処理装置において、前記入力された文書を処理するために必要な知識情報を記憶する知識情報記憶手段と、前記入力された文書に記述されている事象の種類を、前記知識情報記憶手段に記憶されている知識情報を参照して特定する事象特定手段と、前記事象特定手段によって特定された事象に係わる属性の属性値を、前記知識情報記憶手段に記憶されている知識情報を参照して前記文書から抽出する属性値抽出手段と、前記知識情報記憶手段に記憶された知識情報を参照して、前記属性値抽出手段によって抽出された属性値と、実世界の実体との対応付けを行う対応付け手段と、前記対応付け手段によって対応付けがなされた属性値と、前記文書または前記文書の格納場所を特定するための情報とを関連付けて記憶する文書記憶手段と、前記属性値と前記クエリとを参照して、対象となる文書に対して検索処理またはクリッピング処理を行う文書抽出手段と、を有することを特徴とする文書処理装置。

【請求項 2】 前記対応付け手段は、前記属性値の中で固有名であるものに対しては、他の属性値を参照してその固有名が示す実世界の実体を特定するとともに、特定された実体を一意に示す所定の情報を付与し、前記文書抽出手段は、前記対応付け手段によって付与された前記所定の情報を参照して、検索処理またはクリッピング処理を行うことを特徴とする請求項 1 記載の文書処理装置。

【請求項 3】 前記対応付け手段は、前記属性値が参照表現である「同」または「両」を含む場合に、それらの参照表現が参照する属性値を特定することを特徴とする請求項 1 記載の文書処理装置。

【請求項 4】 前記文書に含まれているキーワードの出現頻度を参照して対象とする文書の重要度を算出する重要度算出手段を更に有し、前記重要度算出手段は、前記対応付け手段によって参照先が特定された参照表現もキーワードと同等に処理することを特徴とする請求項 3 記載の文書処理装置。

【請求項 5】 前記属性値のうち、数値化可能なものに関しては、対応する数値に変換して正規化する正規化手段を更に有し、前記文書抽出手段は、前記正規化手段によって正規化された情報を参照して、検索処理またはクリッピング処理を行うことを特徴とする請求項 1 記載の文書処理装置。

【請求項 6】 前記正規化手段によって得られた数値が有する単位を、他の所定の単位に変換する単位変換手段

を更に有することを特徴とする請求項 5 記載の文書処理装置。

【請求項 7】 前記文書を構成する所定の文章の時制を取得する時制取得手段と、前記文書が作成された日時を取得する作成日時取得手段とを更に有し、前記正規化手段は、前記時制取得手段によって取得された文書の時制と、前記作成日時取得手段によって取得された作成日時とを参照して、日時または期間を示す属性値の具体的な値を推定することを特徴とする請求項 5 記載の文書処理装置。

【請求項 8】 前記文書に含まれているキーワードの出現頻度を参照して対象とする文書の重要度を算出する重要度算出手段を更に有し、前記重要度算出手段は、前記正規化手段によって推定された日時または期間を加味して重要度を算出することを特徴とする請求項 7 記載の文書処理装置。

【請求項 9】 前記事象特定手段、前記属性抽出手段、および、対応付け手段は、前記クエリに対しても文書の場合と同様の処理を行い、前記文書抽出手段は、前記対応付け手段によって対応付けがなされた文書の属性値とクエリの属性値とを参照して、検索またはクリッピング処理を行うことを特徴とする請求項 1 記載の文書処理装置。

【請求項 10】 入力された文書に対して所定の処理を施して記憶するとともに、与えられたクエリに対応する文書を、記憶されている文書の中から検索またはクリッピングする処理をコンピュータに実行させるプログラムを記録したコンピュータ読み取り可能な記録媒体において、

コンピュータを、前記入力された文書を処理するために必要な知識情報を記憶する知識情報記憶手段、前記入力された文書に記述されている事象の種類を、前記知識情報記憶手段に記憶されている知識情報を参照して特定する事象特定手段、前記事象特定手段によって特定された事象に係わる属性の属性値を、前記知識情報記憶手段に記憶されている知識情報を参照して前記文書から抽出する属性値抽出手段、

前記知識情報記憶手段に記憶されている知識情報を参照して、前記属性値抽出手段によって抽出された属性値と、実世界の実体との対応付けを行う対応付け手段、前記対応付け手段によって対応付けがなされた属性値と、前記文書または前記文書の格納場所を特定するための情報とを関連付けて記憶する文書記憶手段、前記属性値と前記クエリとを参照して、対象となる文書に対して検索処理またはクリッピング処理を行う文書抽出手段、として機能させるプログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】**【0001】**

【発明の属する技術分野】本発明は入力された文書に対して所定の処理を施して記憶するとともに、与えられたクエリに対応する文書を、記憶されている文書の中から検索またはクリッピングする処理を行う文書処理装置およびそのような処理をコンピュータに実行させるプログラムを記録した記録媒体に関する。

【0002】

【従来の技術】近年、インターネットの普及や全文データベースの増加に伴って、個人が利用可能な情報が飛躍的に増加している。

【0003】このような多量の情報の中から所望の情報を取得する方法としては、例えば、得ようとするデータの特徴を記述した検索式（クエリ）をキーとして、検索処理やクリッピング処理等を行う方法が一般的であった。

【0004】

【発明が解決しようとする課題】しかし、従来の大規模な商用オンラインデータベースや全文検索システムでは、検索式の条件を緩やかにすると検索結果に含まれているノイズ（不要なデータ）が増加し、また、逆に厳しくすると検索洩れが発生するなど、ユーザが望む通りのデータを取得することが困難であるという問題があった。

【0005】即ち、従来の文書フィルタリングにおける文書絞り込み処理や文書検索処理では、クエリと文書の内容の一致度や関連度に基づくランキング検索が行われている程度であるので、文書に含まれている重要情報やユーザの検索意図を十分に反映した文書絞り込みを行うことは困難であった。

【0006】そのため、例えば、組織としての「橋本」が検索したいと思っても、「橋本」という地名が含まれた文書が検索されたりすることが多かった。また、20万円台の新製品について検索したい場合には、「二十万円」「20万円」、「二十一万円」、「二五万円」のように、あらゆる可能性を考慮して生成したクエリを用いる必要があった。

【0007】更に、文書が作成された日時を用いて検索することは可能であっても、文書に含まれている日時情報を活用した検索を行うことができないという問題点があった。

【0008】例えば、以下の文では、同じ「1日」でも示している日が異なる。

(a) A社は、1日、B製品を発売する。

(b) A社は、1日、B製品を発売した。

【0009】ここで、この文が作成された日が1997年2月15日だとすると、(a)の場合では、「1日」は1997年3月1日を指し、また、(b)では1997年2月1日を指すことになる。

【0010】従来の方法においては、文書中の日時情報の属性を認知した上で、このような情報を検索に使用（活用）することが困難であるという問題点があった。本発明はこのような点に鑑みてなされたものであり、ユーザの検索意図を十分に反映した文書検索または文書絞り込みを行うことが可能な文書処理装置を提供することを目的とする。

【0011】また、本発明は、ユーザの検索意図を十分に反映した文書検索処理またはクリッピング処理を行う文書処理を行うプログラムを記録した記録媒体を提供することを目的とする。

【0012】

【課題を解決するための手段】図1は、上記目的を達成する本発明の原理図である。本発明は、入力された文書に対して所定の処理を施して記憶するとともに、与えられたクエリに対応する文書を、記憶されている文書の中から検索またはクリッピングする処理を行う文書処理装置において、知識情報記憶手段3と、事象特定手段4と、属性値抽出手段5と、対応付け手段10と、文書記憶手段11と、文書抽出手段12とから構成されている。

【0013】ここで、知識情報記憶手段3は、入力された文書を処理するために必要な知識情報を記憶する。事象特定手段4は、入力された文書に記述されている事象の種類を、知識情報記憶手段3に記憶されている知識情報を参照して特定する。属性値抽出手段5は、事象特定手段4によって特定された事象に係わる属性の属性値を、知識情報記憶手段3に記憶されている知識情報を参照して文書から抽出する。対応付け手段10は、知識情報記憶手段3に記憶されている知識情報を参照して、属性値抽出手段5によって抽出された属性値と、実世界の实体との対応付けを行う。文書記憶手段11は、対応付け手段10によって対応付けがなされた属性値と、文書または文書の格納場所を特定するための情報とを関連付けて記憶する。文書抽出手段12は、属性値とクエリとを参照して、対象となる文書に対して検索処理またはクリッピング処理を行う。

【0014】知識情報記憶手段3には、事象とそれに関する属性、および、属性を構成する属性値を抽出するための情報とが関連付けられて記憶されている。事象特定手段4は、入力された文書と、知識情報記憶手段3に記憶されている知識情報とを照合することにより、文書中に記述されている事象を特定する。属性値抽出手段5は、知識情報記憶手段3を参照して、特定された事象に関する属性の属性値を文書から抽出する。対応付け手段10は、知識情報記憶手段3に記憶されている知識情報を参照して、抽出された属性値と実世界の实体とを1対1に対応付ける処理を行う。文書記憶手段11は、対応付けがなされた属性値と文書または文書の格納場所を特定するための情報とを関連付けて記憶する。文書抽出手

段 12 は、入力されたクエリに含まれている情報と、文書記憶手段 11 に記憶されている属性値とを照合することにより所望の文書を抽出する。

【0015】これにより、対象となる文書に記述されている内容を、事象という観点から把握し、把握した事象を構成する属性の属性値を抽出するとともに、抽出した属性値を実世界の实体と対応付けることによって生成された情報を参照して文書を検索またはクリッピングするようにしたので、検索またはクリッピングの精度を向上させることが可能となる。

【0016】

【発明の実施の形態】以下、本発明の実施の形態を図面を参照して説明する。図 1 は、本発明の実施の形態の構成例を示すブロック図である。この図において、文書入力部 1 からは、処理対象の文書が入力される。ユーザインタフェース部 2 は、ユーザからのクエリを受け付ける。

【0017】知識情報記憶手段 3 は、後述する事象やその事象に係わる属性に関する情報および固有な固有コードに変換するための情報を記憶している。事象特定手段 4 は、知識情報記憶手段 3 に記憶されている知識情報（事象の種類に関する情報）を参照して、文書入力部 1 またはユーザインタフェース部 2 から入力された文書やクエリに記述されている事象の種類を特定する。

【0018】ここで、「事象」とは、実世界で生起する「できごと」を示している。例えば、新聞記事などでは「A社がXを発売する。」というような実世界で発生した（または、これから発生する）事象に種々の補足情報が付加されて記述されていると考える。

【0019】従って、事象特定手段 4 に対して、例えば、前述の「A社がXを発売する。」が入力されると、この文章に記述されている事象は、＜新製品の発売＞であると特定されることになる。ここで、＜＞は、その内部の語句が抽象化されて得られた概念であることを示している。

【0020】なお、新聞記事などのように、記述の対象となる事象が明確であり、また、その表現様式が限られている文書においては、記述された事象のとりうる構造（以下、事象構造と適宜略記する）にも所定の制約条件が自ずと生ずることから、このような「事象」に着目して文書を解析することにより、効果的な処理を行うことが可能となる。

【0021】属性値抽出手段 5 は、知識情報記憶手段 3 に記憶されている知識情報（所定の事象に係わる属性に関する情報）を参照して、文書またはクエリから属性値を抽出する。

【0022】例えば、知識情報記憶手段 3 は、前述の＜新製品の発売＞という事象に関しては、＜販売会社＞、＜製品情報＞、＜日付＞、および、＜変更点＞などの属性を記憶しており、属性値抽出手段 5 は、事象特定手段

4 によって特定された事象に対応する属性を知識情報記憶手段 3 から取得し、その属性に対応する属性値を文書またはクエリから抽出する。

【0023】例えば、前述の「A社がXを発売する。」という事象では、属性＜販売会社＞に対応する属性値「A社」が取得され、また、属性＜製品情報＞に対応する属性値「X」などが抽出される。

【0024】作成日時取得手段 6 は、文書またはクエリの作成日時を取得する。時制取得手段 7 は、文書またはクエリを構成する文章の時制（tense）を取得する。正規化手段 8 は、属性値抽出手段 5 によって抽出された属性のうち、数値に変換可能なものを選択し、対応する数値に変換（正規化）する。

【0025】単位変換手段 9 は、正規化手段 8 が正規化した数値の単位を変換する処理を行う。対応付け手段 10 は、知識情報記憶手段 3 に記憶されている知識情報を参照して、属性値抽出手段 5 によって抽出された属性値を、実世界における实体に対応付ける処理を行う。なお、この「実体」とは、文書に記述されている属性値が示す実世界の「オブジェクト」を意味している。例えば、前述の例では、「A社」と呼ばれる企業が複数存在する場合には、文書中に記述されている「A社」がどの企業（オブジェクト）を示しているのかを特定する必要があるため、対応付け手段 10 は、文書中の他の属性値（例えば、「社長の名前」や「本社地」など）を参照して「A社」を特定する。

【0026】文書記憶手段 11 は、対応付け手段 10 によって対応付けがなされた属性値集合と、もとの文書（または、もとの文書の格納場所を特定する情報）とを対応付けて記憶する。

【0027】文書抽出手段 12 は、対応付け手段 10 から供給されたクエリに対応する文書を、属性値を参照して文書記憶手段 11 から取得する。そして、重要度算出手段 13 によって算出された個々の文書の重要度を参照し、ある閾値以上の重要度を有する文書を出力する。

【0028】重要度算出手段 13 は、所定のキーワードの出現頻度などを求めることにより、対象とする文書の重要度を算出する。図 2 を参照して、図 1 に示す実施の形態を含む通信システムの構成例について説明する。

【0029】図 2 において、図 1 に示す文書処理装置 20 は、例えば、インターネットなどのネットワーク 21 に接続されている。このネットワーク 21 には、端末装置 22a、22b や、サーバ 23 などが接続されている。

【0030】端末装置 22a、22b は、その入力部をユーザが操作して入力したクエリなどを受け付けて、文書処理装置 20 に送信するとともに、送信したクエリに対応する文書を文書処理装置 20 が送信した場合には、その文書を受信して、例えば、CRT（Cathode Ray Tube）モニタなどへ表示出力する。

10

20

30

40

50

【0031】サーバ23は、記憶部23aに記憶している文書や画像などの情報を、ネットワーク21を介して、要求を行った装置に対して送信する。文書処理装置20は、端末装置22a、22bなどから送信されたクエリを記憶しておき、例えば、サーバ23から新たな文書が供給された場合には、その文書と前述のクエリとの関連度が高い場合には、端末装置22aまたは端末装置22bに対して送信する。

【0032】次に、以上の実施の形態の動作について説明する。図3は、図1に示す実施の形態において、新たな文書が文書入力部1から入力された場合（例えば、図2に示すサーバ23から新たな文書が供給された場合）に実行される処理の一例を説明するフローチャートである。

【0033】このフローチャートが開始されると、以下の処理が実行されることになる。

【S1】文書入力部1は、新たな文書を入力する。

【S2】事象特定手段4は、文書に記述されている事象の種類を特定する。

【0034】即ち、事象特定手段4は、知識情報記憶手段3に記憶されている事象と表現とのマッピング規則情報（図5参照）を参照して、文書に記述されている事象の種類を特定する。図5に示すマッピング規則では、“module” “end” で囲繞された部分が一つの事象（または実体）と表現とのマッピング規則であり、1つの事象に対する表現のバリエーションを記述している。なお、図5に示すマッピング規則の詳細については後述する。

【S3】属性値抽出手段5は、知識情報記憶手段3に記憶されている知識情報を参照して属性値を抽出する。

【0035】例えば、属性値抽出手段5は、図5に示す「module main」内に記述されている事象のバリエーションのうち、入力された文書に対応する定義に含まれている属性（例えば、＜会社情報＞、＜製品＞等）の属性値を、他の「module」や「def」などを参照することにより文書から抽出する。例えば、属性＜会社情報＞に対応する属性値は、第17行目～第21行目に記述されている定義および、第12行目～第15行目に記述されている定義を参照してパターンマッチングを行うことにより、文書から抽出する。

【S4】正規化手段8は、抽出された属性値に日付表現が含まれているか否かを判定し、日付表現が含まれている場合にはステップS5に進み、それ以外の場合にはステップS7に進む。

【S5】作成日時取得手段6は、文書が作成された日時を取得し、また、時制取得手段7は、対象となる事象が記述されている文章の時制（tense）を取得する。

【S6】正規化手段8は、取得した文書作成日時情報と時制情報とを参照して、日付表現を対応する数値に変換する「日付表現変換処理」を実行する。

【0036】なお、この処理の詳細については、図6を参照して後述する。

【S7】正規化手段8は、抽出された属性値に金額表現が含まれているか否かを判定し、金額表現が含まれている場合には、ステップS8に進み、それ以外の場合にはステップS11に進む。

【S8】正規化手段8は、注目している金額表現が規定された通貨単位であるか否かを判定し、規定された通貨単位である場合にはステップS10に進み、それ以外の場合にはステップS9に進む。

【0037】例えば、規定されている通貨単位が「円」ある場合に、「\$」を単位とする金額表現が存在している場合にはステップS9に進む。

【S9】単位変換手段9は、内蔵されている記憶部に記憶している交換レートを読み出し、規定されている通貨単位に変換する処理を行う。

【0038】例えば、「\$100」という表現が存在している場合に、交換レートが「1\$=130円」であるとする、「\$100」は「13000円」に変換されることになる。

【S10】正規化手段8は、金額表現を数値に変換する「金額表現変換処理」を行う。なお、この処理の詳細は、図13を参照して詳述する。

【0039】前述の例では、「13000円」（文字列）が「13000」（数値）に変換される。

【S11】正規化手段8は、他の数値表現が存在するか否かを判定し、他の数値表現が存在する場合にはステップS12に進み、それ以外の場合にはステップS13に進む。

【0040】例えば、「出荷台数50000台」などの表現が存在する場合には、ステップS12に進む。

【S12】正規化手段8は、属性値に含まれている数値表現を対応する数値に変換する。例えば、前述の例では、「50000」（文字列）が計算可能な数値「5000」に変換されることになる。

【S13】対応付け手段10は、属性値に固有名（例えば、「橋本電気」等）が含まれているか否かを判定し、含まれている場合にはステップS14に進み、それ以外の場合にはステップS15に進む。

【S14】対応付け手段10は、固有名を抽出し、その固有名に対応する固有名コードを知識情報記憶手段3から取得して付与する。

【0041】例えば、前述の「橋本電気」に対応する固有名コード「00011」が、知識情報記憶手段3から読み出されて付与される。なお、知識情報記憶手段3には、関連する固有名を相互に関連づけて生成した情報が記憶されているので、文書中のある固有名が複数の候補を有する場合には、関連付けられている他の固有名を参照することにより、固有名を正確に特定することが可能となる。

【0042】即ち、「橋本電気」に対応する候補が「橋本電気株式会社」と「株式会社橋本電気」の2つである場合（同名の会社が存在する場合）には、文書中に記載されている、例えば、社長名や所在地などと、知識情報記憶手段3に関連付けられて記憶されている固有名とを比較することにより、これらの候補を絞り込んで正確な固有名を取得することができる。

【S15】対応付け手段10は、参照表現（「同」または「両」などの表現）が存在するか否かを判定し、参照表現が存在する場合にはステップS16に進み、それ以外の場合にはステップS18に進む。

【0043】例えば、参照表現である（同）が存在する場合には、ステップS16に進むことになる。

【S16】対応付け手段10は、参照表現が参照している対象を同定する。

【0044】例えば、「橋本電気（中山社長）は、橋本電算機（同）の独立を・・・」の場合では、参照表現「同」が参照している対象として「中山社長」を同定する。なお、この同定方法としては、「同」や「両」が内挿された括弧“（）”が検出された場合には、その括弧より前に出現する括弧内の属性値を、参照表現が参照している対象であると同定するようにすればよい。

【S17】対応付け手段10は、参照表現が参照している対象に対応する固有コードを取得し、取得した固有コードを参照表現に対して付与する。

【0045】前述の例では、「中山社長」の固有コード「0001」が参照表現「同」に付与されることになる。

【S18】対応付け手段10は、正規化された属性値（以下、正規化情報と略記する）と、元の文書（または、元の文書の格納場所を特定するための情報）とを関連付けて文書記憶手段11に記憶させる。

【0046】以上の処理により、入力された文書に記述されている事象が特定されるとともに、その事象に係わる属性の属性値が取得される。そして、取得された属性値と実世界の实体とが対応付けられて得られた正規化情報と、もとの文書（または、元の文書の格納場所を特定するための情報）とが文書記憶手段11に記憶されることになる。

【0047】次に、図3のステップS6に示す、「日付表現変換処理」の詳細について説明する。図6は、図3に示す「日付表現変換処理」の詳細を説明するフローチャートである。このフローチャートが開始されると、以下の処理が実行されることになる。

【S30】作成日時取得手段6は、文書の作成日時を取得して、%docyear, %docmonth, %docday に代入する。なお、文書の作成日時は、例えば、新聞記事であればその記事の発行日を取得する。また、新聞記事以外の文書であれば、ファイルの属性を参照して、その作成日時を取得する。

【S31】正規化手段8は、属性値から日付表現を抽出する。

【0048】例えば、対象となる文章が「橋本電気は新たなコンピュータを1日に発売。」であれば、日付表現として「1日」を抽出する。

【S32】正規化手段8は、抽出した日付表現が数字と「年」、「月」、または、「日」によって構成されているか否かを判定し、YESと判定した場合にはステップS33に進み、それ以外の場合にはステップS34に進む。

【0049】例えば、前述の「1日」の例では、数字「1」と「日」とによって構成されていることから、YESと判定されてステップS33に進む。

【S33】正規化手段8は、数字変換テーブル（図7参照）を参照して、日付表現を変換する処理を行う。

【0050】図7に示す数字変換テーブルでは、数字表現に対応する正規化数値が関連付けられており、ある数字表現（文字列）が与えられると、その表現に対応する数値が与えられることになる。

【S34】正規化手段8は、図8に示す日時表現変換テーブルを参照して、日時表現に対応する数値に変換する。

【0051】図8に示す日時表現変換テーブルでは、表現と、そのタイプと、対応する正規化数値とが対応付けられている。なお、タイプとは表現の型であり、例えば、「date」は特定の日時を示し、また、「daterange」は所定の期間を示している。例えば、1998年に作成された文書に「昨年の3月4日」という表現があれば、%year には (%docyear-1) = (1998-1) = 1997 が代入され、また、%monthと%dayには、それぞれ“3”と“4”が代入される。

【0052】また、1997年に作成された文書に「1998年の春」という表現があれば、%year には“1998”が代入されるので、from 1998-3-1 to 1998-5-30という正規化値が得られることになる。

【0053】なお、この日時表現テーブルは、一例であるので、図示したテーブル以外にも種々の実施の形態が考えられる。

【S35】正規化手段8は、全ての値が取得されたか否かを判定し、取得されたと判定した場合にはステップS37に進み、それ以外の場合にはステップS36に進む。

【0054】例えば、年月日に対応する全ての値が取得された場合にはステップS37に進む。

【S36】正規化手段8は、日付の推定処理を行う。なお、この処理の詳細は、図9を参照して後述する。

【S37】正規化手段8は、正規化された数値を%year, %month, %dayにそれぞれ代入して処理を終了する。

【0055】以上の処理によれば、文書に含まれている日付表現を、対応する数値に変換することが可能とな

る。次に、図 9 を参照して図 6 のステップ S 3 6 に示す「日付推定処理」の詳細について説明する。図 9 に示すフローチャートが開始されると、以下の処理が実行されることになる。

〔S 5 0〕正規化手段 8 は、%year のみ値が未代入であるか否かを判定し、未代入である場合にはステップ S 5 1 に進み、それ以外の場合にはステップ S 5 2 に進む。

〔S 5 1〕正規化手段 8 は、%year 推定処理を行う。なお、この処理の詳細については、図 1 0 を参照して後述する。

〔S 5 2〕正規化手段 8 は、%day 以外の値が未代入であるか否かを判定し、未代入である場合にはステップ S 5 3 に進み、それ以外の場合にはステップ S 5 5 に進む。

〔S 5 3〕正規化手段 8 は、%month 推定処理を行う。なお、この処理の詳細については、図 1 1 を参照して詳述する。

〔S 5 4〕正規化手段 8 は、%year 推定処理を行う。

〔S 5 5〕正規化手段 8 は、%month 以外の値が未代入であるか否かを判定し、その結果、未代入である場合にはステップ S 5 6 に進み、それ以外の場合にはステップ S 5 8 に進む。

〔S 5 6〕正規化手段 8 は、%day 推定処理を行う。なお、この処理の詳細については図 1 2 を参照して詳述する。

〔S 5 7〕正規化手段 8 は、%year 推定処理を行う。

〔S 5 8〕正規化手段 8 は、%year 以外の値が未代入であるか否かを判定し、未代入の場合にはステップ S 5 9 に進み、それ以外の場合には処理を終了する。

〔S 5 9〕正規化手段 8 は、推定日時を「from %year-1-1 to %year-12-31」とする。即ち、%year 以外の値が未代入である場合には、検索漏れが生ずることを防止するために、正規化値をできるだけ広い値に設定する。

【0 0 5 6】次に、図 1 0 を参照して、図 8 のステップ S 5 1、S 5 4、S 5 7 に示す「%year 推定処理」の詳細について説明する。このフローチャートが開始されると、以下の処理が実行されることになる。

〔S 6 0〕正規化手段 8 は、対象となる文章から時制取得手段 7 によって取得された時制が未来形である場合にはステップ S 6 1 に進み、それ以外の場合にはステップ S 6 5 に進む。

〔S 6 1〕正規化手段 8 は、%docmonth が %month よりも大きいかなんかを判定し、大きいと判定した場合にはステップ S 6 2 に進み、それ以外の場合にはステップ S 6 3 に進む。

〔S 6 2〕正規化手段 8 は、%year に値 (%docyear+1) を代入する。

【0 0 5 7】例えば、文書が作成された月が 4 月である場合に、「3 月に・・・する予定である」という表現が文章中にある場合には、この「3 月」は、来年の 3 月を示していると推定されることから、%year には値 (%doc

year+1) が代入される。

〔S 6 3〕正規化手段 8 は、%docmonth が %month 以下であるか否かを判定し、YES と判定した場合にはステップ S 6 4 に進み、それ以外の場合にはステップ S 6 5 に進む。

〔S 6 4〕正規化手段 8 は、%year に %docyear の値を代入する。

〔S 6 5〕正規化手段 8 は、時制取得手段 7 によって取得された時制が過去であるか否かを判定し、YES と判定した場合にはステップ S 6 6 に進み、それ以外の場合には図 9 の処理に復帰 (リターン) する。

〔S 6 6〕正規化手段 8 は、%docmonth の値が %month の値以上であるか否かを判定し、YES と判定した場合にはステップ S 6 7 に進み、それ以外の場合にはステップ S 6 8 に進む。

〔S 6 7〕正規化手段 8 は、%docyear の値を %year に代入する。

〔S 6 8〕正規化手段 8 は、%docmonth の値が %month の値よりも小さいかなんかを判定し、YES と判定した場合にはステップ S 6 9 に進み、それ以外の場合には図 9 の処理に復帰する。

〔S 6 9〕正規化手段 8 は、値 (%docyear-1) を %year に代入する。

【0 0 5 8】例えば、文書が作成された月が 4 月である場合に、「6 月に・・・した」という表現が文章中にある場合には、この「6 月」は、昨年の 6 月を示していると推定されることから、%year には値 (%docyear-1) が代入される。

【0 0 5 9】次に、図 1 1 を参照して、図 9 のステップ S 5 3 に示す「%month 推定処理」の詳細について説明する。このフローチャートが開始されると、以下の処理が実行されることになる。

〔S 8 0〕正規化手段 8 は、時制取得手段 7 によって取得された対象となる文章の時制が未来形である場合にはステップ S 8 1 に進み、それ以外の場合にはステップ S 8 5 に進む。

〔S 8 1〕正規化手段 8 は、%docday が %day よりも大きいかなんかを判定し、YES と判定した場合にはステップ S 8 2 に進み、それ以外の場合にはステップ S 8 3 に進む。

〔S 8 2〕正規化手段 8 は、%month に値 (%docmonth+1) を代入する。

【0 0 6 0】例えば、文書が作成された日が 2 日である場合に、「4 日に・・・する予定である」という表現が文章中にある場合には、この「4 日」は、次の月の 4 日を示していると推定されることから、%month には値 (%docmonth+1) が代入されることになる。

〔S 8 3〕正規化手段 8 は、%docday の値が %day の値以下であるか否かを判定し、YES と判定した場合にはステップ S 8 4 に進み、それ以外の場合にはステップ S 8

5に進む。

〔S 8 4〕正規化手段 8 は、%month に %docmonth の値を代入する。

〔S 8 5〕正規化手段 8 は、時制取得手段 7 によって取得された時制が過去であるか否かを判定し、YES と判定した場合にはステップ S 8 6 に進み、それ以外の場合には図 9 の処理に復帰（リターン）する。

〔S 8 6〕正規化手段 8 は、%docday の値が %day の値以上であるか否かを判定し、YES と判定した場合にはステップ S 8 7 に進み、それ以外の場合にはステップ S 8 8 に進む。

〔S 8 7〕正規化手段 8 は、%docmonth の値を %month に代入する。

〔S 8 8〕正規化手段 8 は、%docday の値が %day の値よりも小さいか否かを判定し、YES と判定した場合にはステップ S 8 9 に進み、それ以外の場合には図 9 の処理に復帰する。

〔S 8 9〕正規化手段 8 は、値 (%docmonth-1) を %month に代入する。

【0061】例えば、文書が作成された日が 4 日である場合に、「6 日に・・・した」という表現が文章中にある場合には、この「6 日」は、前の月の 6 日を示していると推定されることから、%month には値 (%docmonth-1) が代入される。

【0062】次に、図 12 を参照して、図 9 のステップ S 5 6 に示す「%day 推定処理」の詳細について説明する。このフローチャートが開始されると、以下の処理が実行されることになる。

〔S 1 0 0〕正規化手段 8 は、%month の値が 1, 3, 5, 6, 8, 10, または、12 のうちの何れかであるか否かを判定し、YES と判定した場合にはステップ S 1 0 1 に進み、それ以外の場合にはステップ S 1 0 2 に進む。

〔S 1 0 1〕正規化手段 8 は、日付情報として「from %year-%month-1 to %year-%month-31」を生成する。

〔S 1 0 2〕正規化手段 8 は、%month の値が 4, 6, 9, または、11 のうちの何れかであるか否かを判定し、YES と判定した場合にはステップ S 1 0 3 に進み、それ以外の場合にはステップ S 1 0 4 に進む。

〔S 1 0 3〕正規化手段 8 は、日付情報として「from %year-%month-1 to %year-%month-30」を生成する。

〔S 1 0 4〕正規化手段 8 は、「年」に関する属性値を参照して、閏年か否かを判定し、閏年である場合にはステップ S 1 0 5 に進み、それ以外の場合にはステップ S 1 0 6 に進む。

〔S 1 0 5〕正規化手段 8 は、日付情報として「from %year-%month-1 to %year-%month-29」を生成する。

〔S 1 0 6〕正規化手段 8 は、日付情報として「from %year-%month-1 to %year-%month-28」を生成する。

【0063】以上の処理によれば、文書に含まれている

日付情報が不十分な情報しか含んでいない場合においても、文書の作成日時と、注目する文章の時制とに応じて日付情報を推定するようにしたので、文書に含まれている日付情報を検索を行う際に有効に活用することが可能となる。

【0064】例えば、「来年の春」のような曖昧な表現も、具体的な数値（例えば、1998 年 3 月 1 日～1998 年 5 月 31 日）に変換（正規化）することが可能となるので、このような曖昧な表現も検索の際に活用することが可能となる。

【0065】次に、図 13 を参照して、図 3 のステップ S 1 0 に示す「金額表現変換処理」の詳細について説明する。このフローチャートが開始されると以下の処理が実行されることになる。

〔S 1 2 0〕正規化手段 8 は、図 14 に示す金額表現変換テーブルを参照して、金額表現を対応する数値に変換し、変数 x に代入する。

【0066】例えば、「二十万円」という表現では、先ず、「二」が“2”に変換され、「十」が“×10”に、また、「万」が“×10000”に変換されるので、全体として値“200000”が得られることになる。

〔S 1 2 1〕正規化手段 8 は、金額表現が「以上」で終わるか否かを判定し、「以上」で終わる場合にはステップ S 1 2 2 に進み、それ以外の場合にはステップ S 1 2 3 に進む。

〔S 1 2 2〕正規化手段 8 は、正規化表現として「from x to *」を生成する。ここで、「*」は任意の値を意味している。

【0067】前述の例では、x=2000 であるので、「from 2000 to *」が生成される。

〔S 1 2 3〕正規化手段 8 は、金額表現が「以下」で終わるか否かを判定し、「以下」で終わる場合にはステップ S 1 2 4 に進み、それ以外の場合にはステップ S 1 2 5 に進む。

〔S 1 2 4〕正規化手段 8 は、正規化表現として「from * to x」を生成する。

〔S 1 2 5〕正規化手段 8 は、金額表現が「台」で終わるか否かを判定し、「台」で終わる場合にはステップ S 1 2 6 に進み、それ以外の場合にはステップ S 1 2 8 に進む。

〔S 1 2 6〕正規化手段 8 は、正規化表現として「from x to x」を生成する。

〔S 1 2 7〕正規化手段 8 は、「to」の後の x に含まれている“0”を“9”に全て変換する。

【0068】例えば、「10 万円台」という表現では、x=100000 となるので、この場合には「to」の後の x に含まれている“0”が“9”に全て変換されるので、199999 となる。従って、正規化表現としては、「from 100000 to 199999」が生成されることになる。

〔S 1 2 8〕正規化手段 8 は、金額表現が「台前半」で

終わるか否かを判定し、「台前半」で終わる場合にはステップ S 1 2 9 に進み、それ以外の場合にはステップ S 1 3 1 に進む。

【S 1 2 9】正規化手段 8 は、正規化表現として「from x to x」を生成する。

【S 1 3 0】正規化手段 8 は、「to」の後の x に含まれている最初の「0」を「5」に変換する。

【0 0 6 9】例えば、「1 0 万円台前半」という表現では、x = 100000 となる。この場合には「to」の後の x に含まれている最初の「0」が「5」に変換されるので、150000 となる。従って、正規化表現としては、「from 1 00000 to 150000」が生成されることになる。

【S 1 3 1】正規化手段 8 は、金額表現が「台後半」で終わるか否かを判定し、「台後半」で終わる場合にはステップ S 1 3 2 に進み、それ以外の場合には図 3 の処理に復帰する。

【S 1 3 2】正規化手段 8 は、正規化表現として「from x to x」を生成する。

【S 1 3 3】正規化手段 8 は、「from」の後の x に含まれている最初の「0」を「6」に変換する。

【S 1 3 4】正規化手段 8 は、「to」の後の x に含まれている「0」を「9」に変換する。

【0 0 7 0】例えば、「1 0 万円台後半」という表現では、x = 100000 となるので、この場合には「to」の後の x に含まれている最初の「0」が「6」にステップ S 1 3 3 において変換され、また、「to」の後の x に含まれている「0」が「9」に全て変換されるので、正規化表現としては、「from 160000 to 199999」が生成されることになる。

【0 0 7 1】以上の処理によれば、例えば、漢数字によって記述されている金額表現を対応する数値に変換するとともに、例えば、「以上」や「台前半」などの曖昧な表現を含む金額表現も対応する数値に変換することが可能となる。

【0 0 7 2】次に、具体的な例を挙げて以上の実施の形態の動作について説明する。いま、図 1 5 に示す文書が図 1 に示す文書入力部 1 から入力されたとする。なお、図 1 5 に示す例文は、新製品の発売に関する文書である。

【0 0 7 3】このような文書が文書入力部 1 から入力されると、事象特定手段 4 は文書に記述されている事象を、知識情報記憶手段 3 に記憶されている知識情報を参照して特定する（図 3 のステップ S 2）。

【0 0 7 4】図 1 5 の例では、図 5 の第 4 行目～1 1 行目に記述されている「module main」の中の第 1 番目の項目（＜会社情報＞は〔、〕？＜日付＞、＜製品＞を発売した。）に該当することから、この文書に記述されている事象が「新製品の発売」とであると判定されることになる。

【0 0 7 5】なお、図 5 に示す知識情報では、事象の定

義が「module main」～「end module」によって囲繞された部分に記述されている。また、事象の定義の中に含まれている、例えば、＜会社情報＞などの属性は、「module」や「def」などにおいて定義されている。例えば、属性＜会社情報＞は、第 1 7 行目～第 2 1 行目の「module」内に定義されており、その内容は、（＜業種＞、＜会社名＞）、（＜業種 2＞&連結語；＜会社名＞）、および、（＜会社名＞）の 3 種類である。

【0 0 7 6】ここで、＜業種＞に関する定義は、第 1 2 行目の「def」の後に記述されており、（. *メーカ | . *会社 | . *大手 | . *開発 | . *販売 | . *製造 | . *業）の中の何れかに該当するものが属性＜業種＞の属性値であるとされる。従って、「パソコンメーカ」や「パソコン大手」などの表現は、＜業種＞の属性値であると判定される。なお、「|」は「or」を意味している。

【0 0 7 7】また、同意語または類義語を含めて定義を行う場合には、第 1 9 行目に示されているように、同意語を含める部分を「&」と「;」の間に挿入する。この例では、「連結語」が同意語または類義語を含む部分となり、その詳細は、第 1 6 行目に定義されており、「連結語」=（を専門とする | である | している | する | の）となる。従って、「オフィスオートメーションを専門とする橋本電機」という表現は、会社情報の第 2 番目の定義（＜業種 2＞&連結語；＜会社名＞）に該当することになる。

【0 0 7 8】このように、本実施の形態においては、トップダウン的な処理が実行されることから、文脈に応じたパターンマッチングが可能となる。以上のような処理によって事象の種類が特定されると、時制取得手段 7 は、事象が記述されている文章を取得し、その時制情報を取得する。図 1 5 に示す文書の例では、その時制は過去形（「発売した」）であるので、「過去形」が時制情報として取得される。なお、このようにして取得された時制は、図 1 6 の第 2 行目に示すように「アスペクト＝過去」として、正規化情報に付加される。

【0 0 7 9】次に、属性値抽出手段 5 は、特定された事象の種類に応じて、属性値を抽出する（図 3 のステップ S 3）。即ち、属性値抽出手段 5 は、図 5 に示す知識情報と文書との間でパターンマッチングを行うことによって属性値を抽出する。

【0 0 8 0】図 1 5 の例では、例えば、＜組織名＞として「橋本電機」が抽出され、また、新たに発売する＜製品情報＞の＜種別＞としては「JCN 互換パソコン」が抽出され、その＜製品名＞としては、「GNW シリーズ」が抽出されている。

【0 0 8 1】続いて、正規化手段 8 は、文書に日付表現が存在するか否かを判定し（図 3 ステップ S 4）、存在する場合には対応する数値に変換する処理を行う。図 1 5 に示す文書では、「十八日」という表現が含まれてい

ることから、正規化手段 8 は、図 3 に示すステップ S 5 において文書作成日時情報と時制情報とを取得して、ステップ S 6 において日付表現変換処理を行う。

【0082】例えば、文書作成日時が「1993 年 10 月 19 日」であるとする、図 16 の第 6 行目に示すように、「発表日付」としてそのタイプが「date」であり、また、その値が「1998-10-18」である情報が正規化情報に付加されることになる。

【0083】続いて、正規化手段 8 は、図 3 のステップ S 7 において、金額表現が存在するか否かを判定する。図 15 に示す文書では、「十七万八千円」などの表現があることから、ステップ S 8 に進み、そこで、規定された通貨単位が否かが判定される。例えば、規定された通貨単位が「円」とし、対象となる表現が前述の「十七万八千円」である場合には、ステップ S 10 に進むことになる。

【0084】なお、「\$ 150」などの表現が含まれている場合には、ステップ S 9 において交換レート（1 \$ *

0001 橋本電機<会社名> 00011 橋本太郎<社長名>
00012 岡山県<所在地>

取得された「橋本電機」に対する候補が複数存在する場合には、橋本電機に関連付けられて記憶されている他の固有名（橋本太郎、岡山県）などが文書中に含まれていないか判定され、候補が絞り込まれることになる。

【0089】そして、ステップ S 14 において、絞り込みの結果得られた固有名コード（例えば、0001）が、正規化情報に付与されることになる（図 16 第 4 行目参照）。

【0090】ステップ S 15 では、対応付け手段 10 は、参照表現が存在するか否かを判定する。図 15 に示す例では、参照表現は存在しないから、NO と判定されてステップ S 18 において、生成された正規化情報と文書（または、文書が格納されている場所を示す情報）とを文書記憶手段 11 に記憶して処理を終了する。

【0091】図 17 は、他の文書例を示している。また、図 18 は、図 17 に示す文書を処理して得られた正規化情報の一例を示している。図 18 の第 3 行目に示すように、図 17 に示す文書に記述されている事象は、合併情報（field = 合併情報）であり、その時制は過去（アスペクト = 過去）である。また、「発表した」という表現が第 1 番目の文章中にないことから、第 2 行目に示すように、「文末表現 = 発表述語なし」とされている。

【0092】更に、第 5 行目から第 7 行目に示されている「合併主体組織情報」の内容としては、第 8 行目と第 18 行目に示されている北海道大木リフトと、東北海道大木リフトとが合併する主体組織であり、それ以外の行には、これらの組織を補足するための＜合併組織補足情報＞が記載されている。

【0093】第 3 4 行目以降には、分析の対象となった

* = 130 円) に応じて、通貨単位の変換が行われた後、ステップ S 10 に進む。

【0085】ステップ S 10 では、文字列「十七万八千円」が、値“178000”に変換される。続くステップ S 11 では、他の数値表現が存在するか否かが判定されるが、図 15 に示す例の第 1 番目の文章には、日付表現以外の数値表現は存在しないことから、ステップ S 13 に進む。

【0086】ステップ S 13 では、対応付け手段 10 が固有名が存在するか否かを判定する。図 15 の例では、固有名「橋本電機」が存在することから、ステップ S 14 に進む。

【0087】ステップ S 14 では、対応付け手段 10 が知識情報記憶手段 3 に記憶されている知識情報のうち、橋本電機に対応する情報を取得する。なお、この情報は、例えば、以下のような情報である。

【0088】

文以外の残りの文章が記載されている。なお、この例では、図 17 の第 3 行目に「同」という参照表現が含まれているので、図 18 の第 2 3 行目に示すように「参照先 = 前」という記述が追加され、参照表現「同」が、第 13 ~ 第 16 行目に示されている「芥川龍太郎 (0251)」（要素 2）であることが示されている。

【0094】次に、以上のようにして生成された正規化情報を参照して、文書を検索する場合の処理の一例について説明する。図 19 は、図 1 に示すユーザインタフェース部 2 に表示される入力画面の表示例である。この例では、製品の販売情報が記載された文書を検索の対象としている。即ち、＜製品の販売＞が事象として記述された文書が検索の対象とされる。

【0095】この例では、第 1 番目に示すボックス「組織名」に、製品を発売した組織名が入力される。また、第 2 番目に示すボックス「製品種」には、製品の種類が入力される。更に、ボックス「価格」には製品の価格の範囲が入力される。ボックス「発売日」には、発売された日の範囲が入力される。なお、最下行に表示されているボタン「検索」は、全ての入力終了し、検索を開始する場合に操作される。

【0096】図 20 は、図 19 に示す画面に所定のクエリが入力された場合の入力例を示している。この例では、組織名として「AAA」が、また、製品種として「パソコン」が入力されている。

【0097】更に、価格は、「100000」円以上「300000」円以下とされており、発売日は「1997」年「1」月「1」日から「1997」年「6」月「30」日までとされている。

【0098】このような入力画面から入力されたクエリ

は、各入力項目の属性を示す情報が付与された後、事象特定手段 4、属性値抽出手段 5、および、対応付け手段 10 を介して、文書抽出手段 12 に供給される。なお、付与される情報としては、例えば、「AAA」に対してはタグ<組織名>が付与される。また、価格はタグ<価格 type=price unit= 円 value = “ from 100000 to 300000 ”>に変換される。更に、発売日は、タグ<発売日 type=date value= “from 1997-1-1 to 1997-6-30 ”>に変換される。

【0099】文書抽出手段 12 は、ユーザインタフェース部 2 から供給されたクエリとタグとに対応する属性値を有する文書を文書記憶手段 11 から取得する。即ち、文書記憶手段 11 には、元の文書とともに正規化情報が記憶されているので、文書抽出手段 12 は、この正規化情報に含まれている属性値と、クエリのタグとを照合することにより、所望の文書を抽出する。

【0100】このようにして検索された結果は、図示せぬ表示装置に表示出力される。図 21 は、検索結果を表示する画面のテンプレートを示している。この例では、検索結果の属性値として「組織名」、「製品種」、「製品名」、「価格」、「発売日」、および、「見出し」が表示される。

【0101】図 22 は、実際の表示例を示している。この例の第 1 行目の項目は、「AAA」という組織が、デスクトップ型のパソコンを、200000~299999 円で、1997/02/29 に発売しており、その文書の見出しは「低価格パソコン発売」であることを示している。

【0102】図 23 は、図 1 に示すユーザインタフェース部 2 に表示される入力画面の他の表示例である。この例では、「組織の合併情報」が記載された文書を検索の対象としている。即ち、組織の合併が事象として記述された文書が検索の対象とされる。この例では、第 1 番目と第 2 番目に示すボックス「組織名」に、合併する組織名が入力される。また、ボックス「合併日」には、合併が行われる日の範囲が入力される。なお、最下行に表示されているボタン「検索」は、全ての入力終了した後、検索を開始する場合に操作される。

【0103】図 24 は、図 23 に示す入力画面に所定のクエリが入力された場合の入力例を示している。この例では、組織名として「AAA」が、また、合併日として「1997」年「1」月「1」日から、「1997」年「12」月「31」日までが入力されている。

【0104】このような入力画面において、ボタン「検索」が操作されると、前述の場合と同様にタグが生成され、文書記憶手段 11 に記憶されている正規化情報と、このタグとを照合することにより、文書が検索される。

【0105】図 25 は、図 24 の検索結果を表示する画面の表示例である。この表示例では、検索結果の属性として「組織名」、「組織名」、「新組織名」、「合併日」、および、「見出し」が表示される。

【0106】図 26 は、実際の表示例を示している。この例では、検索結果の文書には、組織名が「AAA」および「BBB」である会社が「1997/04/01」に合併し、新組織名は「CCC」であることが示されており、また、その文書の見出しは、「AAA, BBB, 2 社合併」であることが示されている。

【0107】以上の実施の形態によれば、検索の対象となる事象に対応した入力画面を用意して、その入力画面から必要な項目を入力することにより、所望の文書が取得されることになる。ところで、文書記憶手段 11 に記憶されている文書には、前述の正規化情報が関連付けられて記憶されているので、その正規化情報を参照することにより、例えば、対象とする文書に、新たに発売されたパソコンの価格が「二十五万円」と漢数字で記載されているような場合においても、「200000」円~「300000」円と記述されたクエリによって取得されることになる。

【0108】なお、以上の実施の形態においては、検索しようとする事象に対応した入力画面から所定の項目を入力し、入力された項目に対応する文書を検索するようにしたが、クエリを文章として入力し、入力された文章に対して正規化処理を行った後、対応する文書を検索するようにしてもよい。以下、そのような方法により、クエリを正規化する処理の一例について、図 27 を参照して説明する。このフローチャートが開始されると以下の処理が実行されることになる。

【S151】ユーザインタフェース部 2 は、文章として記述されたクエリを入力する。

【S152】事象特定手段 4 は、クエリに記述されている事象の種類を特定する。即ち、事象特定手段 4 は、知識情報記憶手段 3 に記憶されている事象と表現とのマッピング規則情報（図 5 参照）を参照して、クエリに記述されている事象の種類を特定する。

【S153】属性値抽出手段 5 は、知識情報記憶手段 3 に記憶されている知識情報を参照して属性値を抽出する。

【S154】正規化手段 8 は、抽出された属性値に日付表現が含まれているか否かを判定し、日付表現が含まれている場合にはステップ S155 に進み、それ以外の場合にはステップ S157 に進む。

【S155】作成日時取得手段 6 は、クエリが作成された日時を取得し、また、時制取得手段 7 は、クエリの時制 (tense) を取得する。

【S156】正規化手段 8 は、取得したクエリ作成日時情報と時制情報とを参照して、日付表現を対応する数値に変換する「日付表現変換処理」を実行する。なお、この処理の詳細については、図 6 を参照して既述したので、その説明は省略する。

【S157】正規化手段 8 は、抽出された属性値に金額表現が含まれているか否かを判定し、金額表現が含まれ

ている場合には、ステップS158に進み、それ以外の場合にはステップS161に進む。

〔S158〕正規化手段8は、注目している金額表現が規定された通貨単位であるか否かを判定し、規定された通貨単位である場合にはステップS160に進み、それ以外の場合にはステップS159に進む。例えば、規定されている通貨単位が「円」ある場合に、「\$」を単位とする金額表現が既述されている場合にはステップS159に進む。

〔S159〕単位変換手段9は、内蔵されている記憶部に記憶している交換レートを読み出し、規定されている通貨単位に変換する処理を行う。

〔O109〕例えば、「\$100」という表現が存在している場合に、交換レートが「1\$=130円」であるとする、「\$100」は「13000円」に変換されることになる。

〔S160〕正規化手段8は、金額表現を数値に変換する「金額表現変換処理」を行う。なお、この処理の詳細は、図13を参照して既述したので、その説明は省略する。

〔O110〕前述の例では、「13000円」（文字列）が「13000」（数値）に変換されることになる。

〔S161〕正規化手段8は、他の数値表現が存在するか否かを判定し、他の数値表現が存在する場合にはステップS162に進み、その他の場合にはステップS163に進む。

〔O111〕例えば、「出荷台数50000台」などが存在する場合には、ステップS162に進む。

〔S162〕正規化手段8は、属性値に含まれている数値表現を対応する数値に変換する。例えば、前述の例では、文字列「50000」が計算可能な数値「5000」に変換されることになる。

〔S163〕対応付け手段10は、属性値に固有名（例えば、「橋本電気」等）が含まれているか否かを判定し、含まれている場合にはステップS164に進み、それ以外の場合にはステップS165に進む。

〔S164〕対応付け手段10は、固有名を抽出し、その固有名に対応する固有コードを知識情報記憶手段3から取得して属性値に付与する。

〔O112〕例えば、前述の「橋本電気」に対応する固有コード「00011」が、知識情報記憶手段3から読み出されて付与される。なお、知識情報記憶手段3には、関連する固有名を相互に関連づけて生成した情報が記憶されているので、ある固有名が複数の候補を有する場合には、関連付けられている他の固有名を参照することにより、固有名を正確に特定することが可能となる。

〔O113〕即ち、「橋本電気」に対応する候補が「橋本電気株式会社」と「株式会社橋本電気」の2つである場合（同名の会社が存在する場合）には、クエリ中に記載されている、例えば、社長名や所在地など、知識情

報記憶手段3に関連付けられて記憶されている固有名とを比較することにより、これらの候補を絞り込んで正確な固有名コードを取得することができる。

〔S165〕対応付け手段10は、参照表現（同または両などの表現）が存在するか否かを判定し、参照表現が存在する場合にはステップS166に進み、それ以外の場合にはステップS168に進む。

〔O114〕例えば、参照表現である（同）が存在する場合には、ステップS166に進むことになる。

〔S166〕対応付け手段10は、参照表現が参照している対象を同定する。

〔O115〕例えば、「橋本電気（中山社長）」は、橋本電算機（同）の独立を・・・」の場合では、参照表現「同」が参照している対象として「中山社長」を同定する。なお、この同定方法としては、「同」や「両」が内挿された括弧「（）」が検出された場合には、その括弧より前に出現する括弧内の属性値を、参照表現が参照している対象であると同定するようにすればよい。

〔S167〕対応付け手段10は、参照表現が参照している対象に対応する固有コードを取得し、取得した固有コードを参照表現に対して付与する。

〔O116〕前述の例では、「中山社長」の固有コード「00010」が参照表現「同」に付与されることになる。

〔S168〕対応付け手段10は、以上のようにして生成されたクエリの正規化情報を、文書抽出手段12に供給する。その結果、文書抽出手段12は、以上のようにして生成されたクエリの正規化情報を参照して、文書記憶手段11に記憶されている文書を検索する。

〔O117〕例えば、クエリとして「橋本酒造が純米酒、橋本を発売した。」が入力された場合には、事象特定手段4は、知識情報記憶手段3に記憶されている知識情報を参照し、入力されたクエリが「新製品の発売」という事象を示していることを特定する。

〔O118〕属性値抽出手段5は、＜組織名＞として「橋本酒造」を抽出し、また、＜製品種＞として「純米酒」を、また、＜製品名＞として「橋本」を抽出する。対応付け手段10は、「橋本酒造」に対応する固有コードが存在する場合には、その固有コードを取得して付与する。例えば、「橋本酒造」の固有コードが「0111」である場合には、＜組織名＞橋本酒造（0111）＜／組織名＞というタグが生成されることになる。

〔O119〕文書抽出手段12は、以上のようにして生成された正規化情報を参照して、文書記憶手段11から該当する文書を抽出する。即ち、文書抽出手段12は、組織名タグと固有コード（0111）が付けられた「橋本酒造」、製品種タグが付けられた「橋本」、および、製品名タグが付けられた「橋本」を含み、その事象が「新製品の発売」である文書を文書記憶手段11から抽出する。

【0120】このような処理によれば、例えば、「橋本さんは、橋本酒造製の純米酒を注文した。」という一文が含まれている文書が検索結果として取得されることを防止することができる。即ち、クエリと文書の正規化情報には、抽出された属性を示すタグがそれぞれ付与されていることから、例えば、＜製品名＞である「橋本」を、＜人名＞と混同すること防止することができる。

【0121】次に、以上の実施の形態により文書をクリッピングする場合の処理の一例について説明する。図29は、文書をクリッピングする場合に、ユーザから送信されたクエリを正規化する処理の一例を説明するフローチャートである。このフローチャートが開始されると、以下の処理が実行されることになる。

【S180】ユーザインタフェース部2は、所定のユーザからのクエリを入力する。【S181】事象特定手段4、属性値抽出手段5、および、対応付け手段10は、図27および図28のステップS151～S167の処理を実行して、クエリを正規化する。

【S182】文書抽出手段12は、正規化されたクエリ（正規化情報）と、それを送信したユーザを特定する情報とを関連付けて記憶する。

【S183】文書抽出手段12と重要度算出手段13は、文書記憶手段11に記憶されている文書と、ユーザ毎のクエリの関連度を判定する「関連度判定処理」を実行する。なお、この処理の詳細は、図30を参照して後述する。

【0122】次に、図30を参照して、図29に示す「関連度判定処理」の詳細について説明する。このフローチャートが開始されると、以下の処理が実行されることになる。

【S201】重要度算出手段13は、正規化情報が付与された文書と、正規化されたクエリとの関連度をユーザ単位で計算する。

【0123】なお、計算方法としては、例えば、正規化されたクエリに含まれている重要表現を対象となる文書がいくつ含んでいるかに応じてスコアリングを行い、スコアの高い文書を関連度の高い文書とする方法を採用することができる。

【S202】文書抽出手段12は、重要度算出手段13の計算結果を参照し、関連度の高い文書を抽出する。

【S203】文書抽出手段12は、正規化したクエリに含まれている日付、金額、および、数値が、文書の正規化情報に含まれているそれらの値と一致する文書を抽出する。

【S204】文書抽出手段12は、一致した文書をネットワーク21を介してユーザに送付する。

【0124】続いて、図31を参照して、例えば、サーバ23から新たな文書が送信されてきた場合に、文書処理装置20において実行される処理の一例を説明する。このフローチャートが開始されると、以下の処理が実行

されることになる。

【S230】文書入力部1は、ネットワーク21を介して、例えば、サーバ23から新たな文書の入力を受ける。

【S231】事象特定手段4、属性値抽出手段5、および、対応付け手段10は、文書の正規化処理を実行する。

【0125】即ち、事象特定手段4、属性値抽出手段5、および、対応付け手段10は、図3および図4に示す処理を実行することにより入力された文書に対応する正規化情報を生成する。

【S232】文書抽出手段12および重要度算出手段13は、図30に示す「関連度判定処理」を実行する。その結果、生成された正規化情報に一致するクエリが存在している場合には、そのクエリを送信したユーザに対して、新たに入力された文書が送付される。

【0126】以上の処理によれば、新たな文書が入力された場合には、入力された文書の正規化情報と、各ユーザの正規化されたクエリとの関連度を算出して、関連度が高い場合には対応するユーザに対して文書を送信するようにしたので、ユーザの要求に適合した文書を正確に選択して送信することが可能となる。

【0127】なお、上記の処理機能は、コンピュータによって実現することができる。その場合、文書処理装置が有すべき機能の処理内容は、コンピュータで読み取り可能な記録媒体に記録されたプログラムに記述されており、このプログラムをコンピュータで実行することにより、上記処理がコンピュータで実現される。コンピュータで読み取り可能な記録媒体としては、磁気記録装置や半導体メモリ等がある。

【0128】市場に流通させる場合には、CD-ROM (Compact Disk Read Only Memory) やフロッピーディスク等の可搬型記録媒体にプログラムを格納して流通させたり、ネットワークを介して接続されたコンピュータの記憶装置に格納しておき、ネットワークを通じて他のコンピュータに転送することもできる。コンピュータで実行する際には、コンピュータ内のハードディスク装置等にプログラムを格納しておき、メインメモリにロードして実行するようにすればよい。

【0129】

【発明の効果】以上説明したように本発明では、対象となる文書に記述されている事象を特定し、特定された事象に関する属性の属性値を抽出し、抽出した属性値と実世界の实体とを対応付けすることによって生成された情報を参照して、文書を検索またはクリッピングするようにしたので、各属性値を正確に認識して文書を検索またはクリッピングすることが可能となるので、結果として、文書の検索またはクリッピング精度を向上させることが可能となる。

【図面の簡単な説明】

【図 1】本発明の実施の形態の構成例を示すブロック図である。

【図 2】図 1 に示す文書処理装置を含む通信システムの構成例である。

【図 3】文書の正規化処理の一例を説明するフローチャートである。

【図 4】文書の正規化処理の一例を説明するフローチャートである。

【図 5】知識情報の一例を示す図である。

【図 6】図 3 に示す日付表現変換処理の詳細を説明するフローチャートである。

【図 7】数字変換テーブルの一例を示す図である。

【図 8】日時表現変換テーブルの一例を示す図である。

【図 9】図 6 に示す日付推定処理の詳細を示すフローチャートである。

【図 10】図 9 に示す%year 推定処理の詳細を説明するフローチャートである。

【図 11】図 9 に示す%month推定処理の詳細を説明するフローチャートである。

【図 12】図 9 に示す%day推定処理の詳細を説明するフローチャートである。

【図 13】図 3 に示す金額表現変換処理の詳細を説明するフローチャートである。

【図 14】金額表現変換テーブルの一例を示す図である。

【図 15】図 1 に示す実施の形態に入力される文書の一例である。

【図 16】図 15 に示す文書を処理した結果生成される正規化情報の一例である。

【図 17】図 1 に示す実施の形態に入力される文書の他の一例を示す図である。

【図 18】図 17 に示す文書を処理した結果生成される正規化情報の一例である。

【図 19】製品販売情報に関する文書を検索する際の入力画面の一例である。

【図 20】図 19 に示す入力画面に入力がなされた場合の一例である。

【図 21】図 19 に示す入力画面に対応する検索結果表*

【図 15】

橋本電機は十八日、「JCN互換パソコン「GNWシリーズ」を発売したと発表した。これまでは独自設計のパソコンを開発、販売してきたが、昨年から低価格競争に対抗するため、世界的な業界標準である「JCN互換機」を投入する。橋本電機が発売するのはデスクトップ型とノート型の合計六機種十七モデル。全機種にジョウフムのCPU（中央演算処理装置）「597」とソフィスディケータのソフトのOS（基本ソフト）、V-OS98を搭載した。個人ユーザーから企業のネットワーク向けまでを対象としている。台湾メーカーから部品供給を受けるなど海外部品調達をこれまでの平均三〇%から七〇%に高め、生産コストを引き下げた。価格はデスクトップ型が十七万八千円から三十二万八千円まで、ノート型はモノクロディスプレイタイプが二十二万八千円から三十四万八千円まで、カラータイプが四十二万八千円から七十四万八千円まで。

* 示画面の一例である。

【図 22】図 20 に示す入力画面に対応する検索結果の画面の一例である。

【図 23】組織合併情報に関する文書を検索する際の入力画面の一例である。

【図 24】図 23 に示す入力画面に入力がなされた場合の一例である。

【図 25】図 23 に示す入力画面に対応する検索結果表示画面の一例である。

【図 26】図 24 に示す入力画面に対応する検索結果の画面の一例である。

【図 27】クエリに対する正規化処理の一例を説明するフローチャートである。

【図 28】クエリに対する正規化処理の一例を説明するフローチャートである。

【図 29】文書のクリッピングを行う場合において、ユーザからのクエリに対する処理の一例を説明するフローチャートである。

【図 30】図 29 に示す関連度判定処理の詳細を説明するフローチャートである。

【図 31】文書のクリッピングを行う場合において実行される、文書に対する処理の一例を説明するフローチャートである。

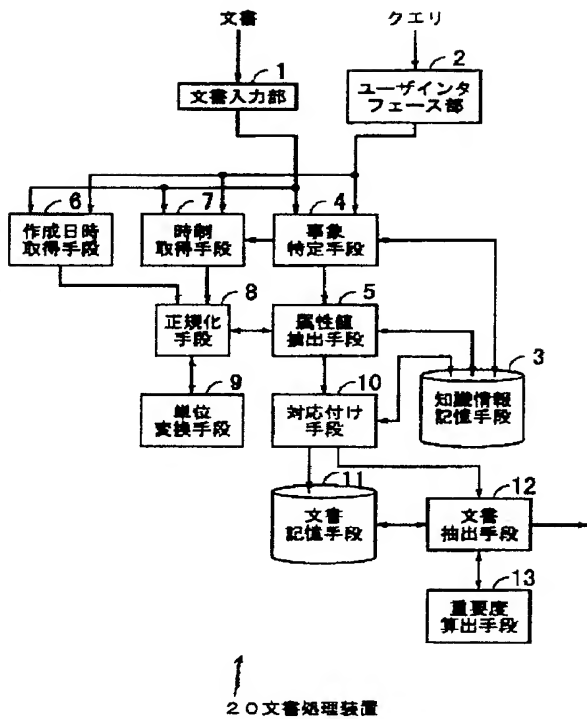
【符号の説明】

- 1 文書入力部
- 2 ユーザインタフェース部
- 3 知識情報記憶手段
- 4 事象特定手段
- 5 属性値抽出手段
- 6 作成日時取得手段
- 7 時制取得手段
- 8 正規化手段
- 9 単位変換手段
- 10 対応付け手段
- 11 文書記憶手段
- 12 文書抽出手段
- 13 重要度算出手段

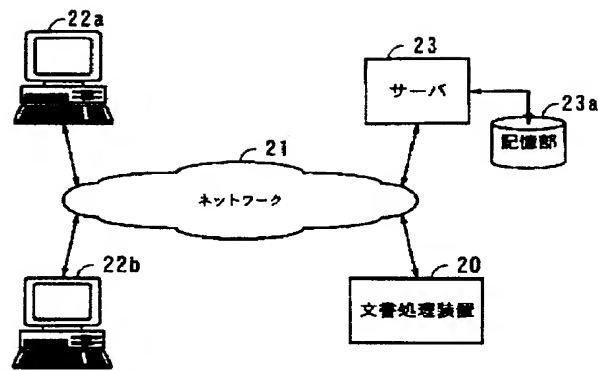
【図 17】

大木リフトの道内販売子会社、北海道大木フォークリフト（石狩管内石狩市、芥川龍太郎社長）と東北海道大木フォークリフト（十勝管内芽室町、同）が一日、合併した。大木リフトが全国で進めている販売強化策に基づくもので、経営を一元化することにより、効率的な営業体制を整える狙い。新会社の名称は北海道大木フォークリフトで、社長には芥川社長が就任した。資本金は二億四千万円、従業員は百三十人。

【図 1】



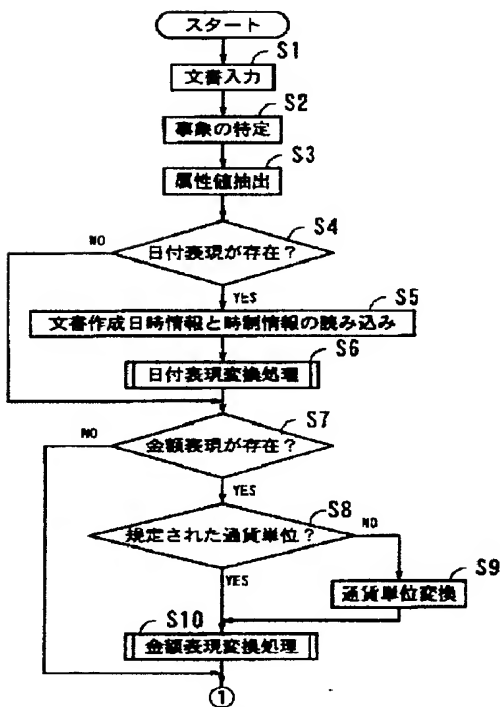
【図 2】



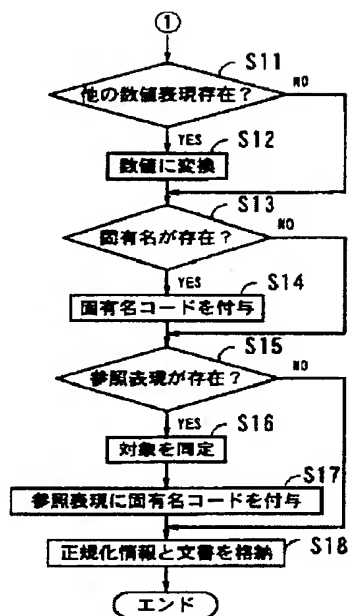
【図 7】

表現	正規化数値	表現	正規化数値
0	0	7	7
〇	0	七	7
1	1	8	8
一	1	八	8
2	2	9	9
二	2	九	9
3	3	10	10
三	3	十	10
4	4	11	11
四	4	十一	11
5	5	12	12
五	5	十二	12
6	6		
六	6		

【図 3】



【図 4】



【図 19】

クエリ入力画面

製品販売情報

組織名

製品種

価格 円以上
 円以下

発売日 年 月 日から
 年 月 日まで

【図5】

```

typedef 名前(.*);
tagclass 名前<会社情報>, <製品>, <変更点>, <行為>
def <日付> ([一二三四五六七八九十]*日)
module main
  <会社情報>は[.]*?<日付>
  <会社情報>は[.]*?<製品>を発売した。
  <会社情報>は[.]*?<製品>を発売した。
  <会社情報>は[.]*?<製品>を[.]?<変更点>して[.]<日付>発売した。
  <会社情報>は[.]*?<製品>を<日付>発売した。
  <会社情報>は[.]*?<製品>を<行為>, <日付>発売した。
endmodule

def<会社名> ([.*]*会社|. *大手|. *開発|. *販売|. *製造|. *業)
def<業種2> ([.*]*)
def<会社名2> ([.*]*)
def<会社補足情報> ([.*]*)
synset 連結語を専門とするである。している。するの

module 会社情報
  <業種>, <会社名>
  <業種2>, <会社名2>
  <会社名2> <連結語> <会社名>
endmodule

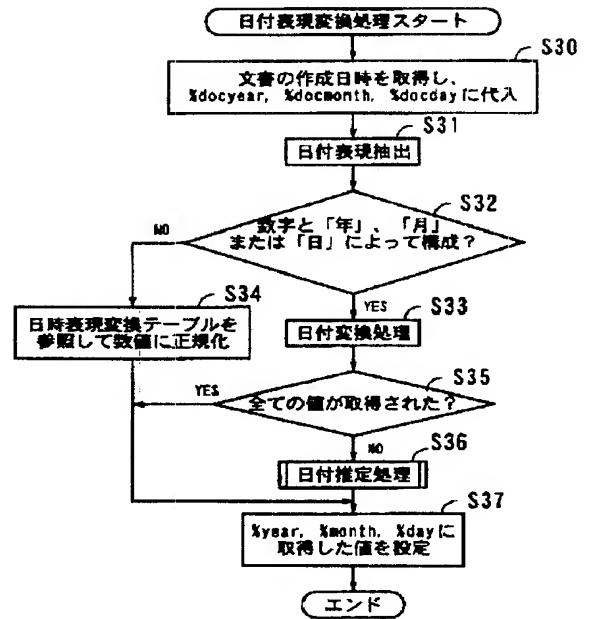
module 会社名
  <会社名> (<会社補足情報>)
endmodule

```

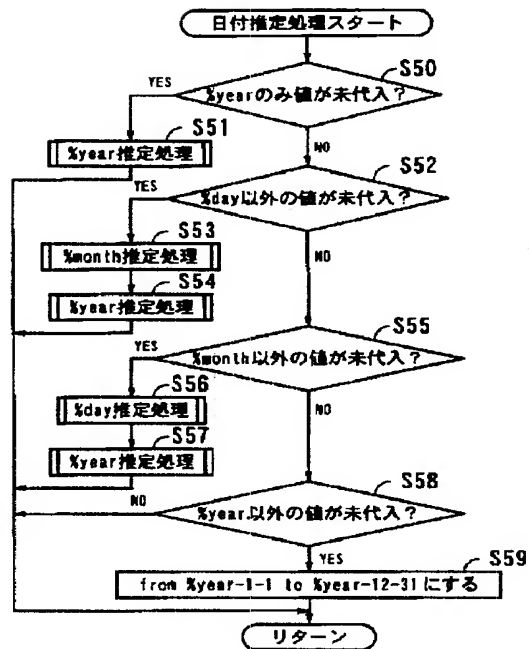
【図8】

表現	type	対応する正規化数値
昨年	date	%docyear - 1
去年	date	%docyear - 1
一昨年	date	%docyear - 2
前年	date	%docyear - 1
来年	date	%docyear + 1
翌年	date	%docyear + 1
本年	date	%docmonth
先月	date	%docmonth - 1
前月	date	%docmonth - 1
来月	date	%docmonth + 1
今月	date	%docmonth
昨日	date	%docday - 1
前日	date	%docday - 1
明日	date	%docday + 1
今日	date	%docday
本日	date	%docday
月初め	daterange	from %year-%month-1 to %year-%month-10
月末	daterange	from %year-%month-21 to %year-%month-31
上旬	daterange	from %year-%month-1 to %year-%month-10
中旬	daterange	from %year-%month-11 to %year-%month-20
下旬	daterange	from %year-%month-21 to %year-%month-31
平成	date	%docyear + 1988
春	daterange	from %year-3-1 to %year-5-31
夏	daterange	from %year-6-1 to %year-8-31
秋	daterange	from %year-9-1 to %year-11-30
冬	daterange	from %year-12-1 to %year-2-29
お正月	daterange	from %year-1-1 to %year-1-7
ゴールデンウィーク	daterange	from %year-4-26 to %year-5-5
GW	daterange	from %year-4-26 to %year-5-5
クリスマス	daterange	from %year-12-1 to %year-12-25
入学シーズン	daterange	from %year-4-1 to %year-4-30
受験シーズン	daterange	from %year-1-10 to %year-1-30
梅雨	daterange	from %year-6-1 to %year-7-10

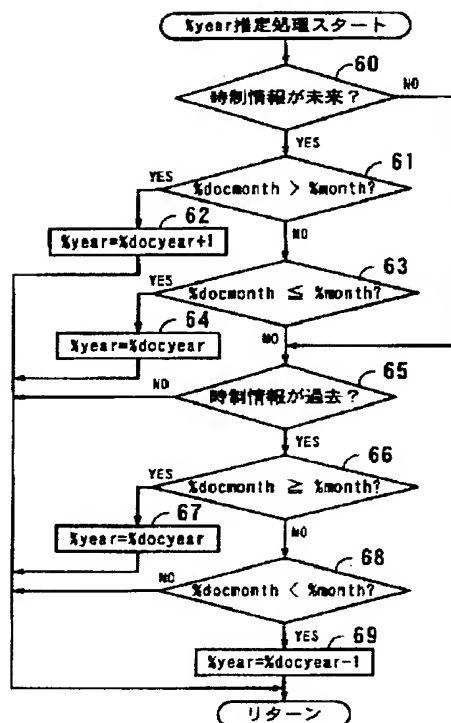
【図6】



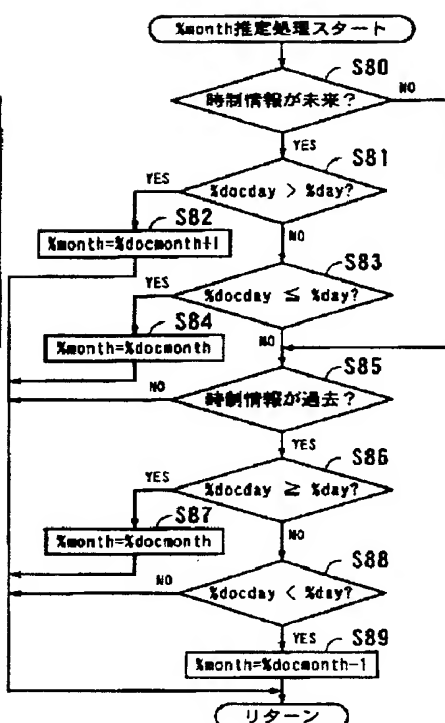
【図9】



【図 10】



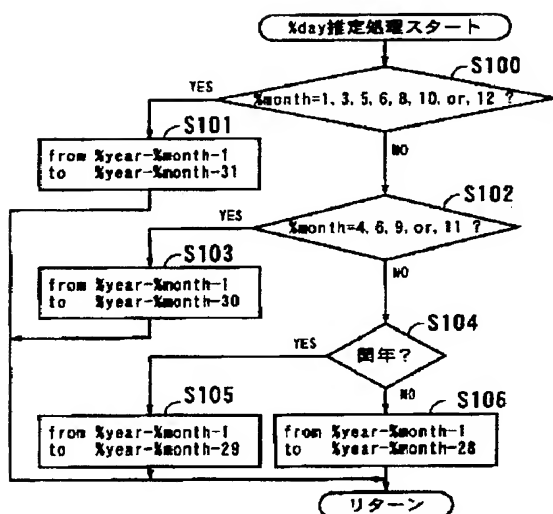
【図 11】



【図 14】

表現	正規化数値	表現	正規化数値
0	0	6	6
〇	0	六	6
ゼロ	0	7	7
零	0	七	7
1	1	8	8
一	1	八	8
壹	1	9	9
2	2	九	9
二	2	十	×10
式	2	百	×100
3	3	千	×1000
三	3	万	×10000
参	3	億	×100000000
4	4	兆	×1000000000000
四	4	京	×10000000000000000
5	5	分	×0.01
五	5	厘	×0.001

【図 12】

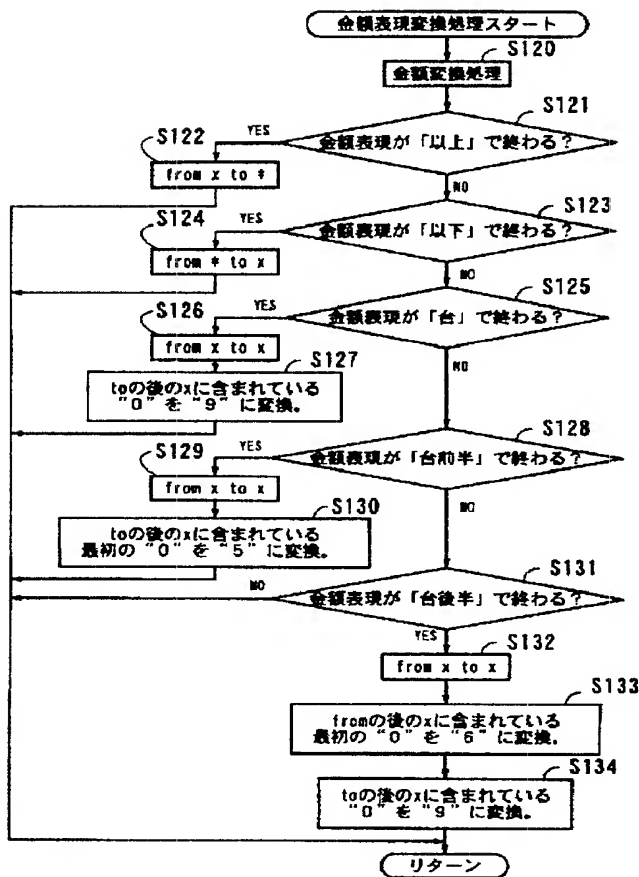


【図 16】

```

<add-anno>
<記事内容 文章表現=発表述語あり アスペクト=過去>
<発表主体組織情報>
<組織名>根本電機(0001)</組織名>
</発表主体組織情報>
<発表日付 type=date value=1993-10-18>
<日>十八日</日>
</発表日付>
<記事内容-用言型1 field=製品販売情報 アスペクト=現在>
<販売情報-用言型1>
<販売情報-用言型10>
<製品情報>
<製品>JCN互換パソコン</製品>
<製品名>GNWシリーズ</製品名>
</製品情報>
</販売情報-用言型10>
</販売情報-用言型1>
</記事内容-用言型1>
</記事内容>
<記事内容残り>これまでは独自設計のパソコンを開発、販売してきたが、
昨年からの低価格競争に対抗するため、世界的な業界標準
であるJCN互換機を導入する。根本電機が販売するのは
デスクトップ型とノート型の合計六機種十七モデル。全機
型にジョブ用のCPU(中央演算処理装置)「697」と
ソフィスティケートソフトのOS(基本ソフト)
V-OSS98を搭載した。個人ユーザーから企業のネット
ワーク向けまでを対象としている。台湾メーカーから部品供給
を受けるなど海外部品調達をこれまでの平均三〇%から七〇%
に高め、生産コストを引き下げた。価格はデスクトップ型が
</記事内容残り>
<価格 type=price unit=円 value="from 178000 to 178000">
<金額>十七万八千</金額>
</価格>
<記事内容残り>から三十七万八千円まで。ノート型はモノクロディスプレ
イタイプが二十二万八千円
  
```

【図13】



【図18】

```

<add-anno>
<記事内容・文本文書=発表述語なし>
<記事内容・用言型0 field=合併情報 アスペクト=過去>
<合併情報・用言型0>
<合併主体組織情報 type=合併主体組織情報>
<組織名>大木リフト020</組織名>
<組織補足情報1 type=組織補足情報>
<組織名>北海道大木リフト021</組織名>
<合併組織補足情報1 type=組織補足情報>
<要項1 要項数=2 type=要項>
<組織所在地>石狩管内石狩市</組織所在地>
</要項1>
<要項2 要項数=2 type=要項>
<氏名>芥川太蔵0251</氏名>
<役職名>社長</役職名>
</要項2>
</合併組織補足情報1>
<組織名>北海道大木リフト021</組織名>
<合併組織補足情報2 type=組織補足情報>
<要項1 要項数=2 type=要項>
<組織所在地>石狩管内石狩市</組織所在地>
</要項1>
<要項2 要項数=2 type=要項 参照タイプ=同1 参照先=前(0251)>
<参照先>同</参照先>
</要項2>
</合併組織補足情報2>
<合併主体組織情報>
<合併日付 type=date value=1994-04-01>
<日>一日</日>
</合併日付>
</記事内容・用言型0>
</記事内容>
</記事内容残り>
<記事内容残り0>大木リフトが全国で最も信頼性の高い販売強化策に基づき、
  顧客を一元化することにより、効率的な営業体
  制を築く。</記事内容残り0>
<新組織名>北海道大木リフト</新組織名>
<記事内容残り0>、社長には芥川社長が就任した。資本金は</記事内容残り0>
</記事内容残り0>
<価格 type=price unit=円 value=from 240000000 to 240000000>
<資本金>二億四千万</資本金>
</価格>
<記事内容残り0>、従業員は百三十人。</記事内容残り0>
</add-anno>

```

【図20】

クエリ入力画面

製品販売情報

組織名 AAA

製品種 パソコン

価格 1000000円以上

3000000円以下

発売日 1997年1月1日から

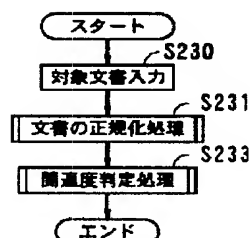
1997年6月30日まで

検索

【図21】

検索結果表示画面					
製品販売情報					
組織名	製品種	製品名	価格	発売日	見出し

【図31】



【図22】

【図23】

検索結果表示画面					
製品販売情報					
組織名	製品種	製品名	価格	発売日	見出し
AAA	デスクトップ型パソコン		200,000~299,999	1997/02/29	低価格パソコン発売
AAA	ディスプレイ一体型パソコン	BBBBB	268,000	1997/03/01	BBBBB新シリーズ発売
AAA	カラーノートパソコン	CCCCC	298,000	1997/04/11	AAA新発売
AAA	互換型パソコン	DDDDD	178,000	1997/05/20	互換型機発売

クエリ入力画面	
組織合併情報	
組織名	<input type="text"/>
組織名	<input type="text"/>
合併日	<input type="text"/> 年 <input type="text"/> 月 <input type="text"/> 日から <input type="text"/> 年 <input type="text"/> 月 <input type="text"/> 日まで
検索	

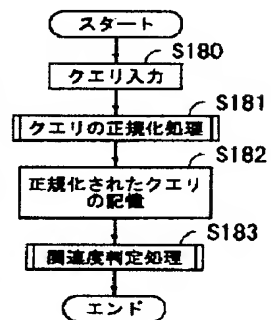
【図24】

【図25】

【図29】

クエリ入力画面	
組織合併情報	
組織名	<input type="text" value="AAA"/>
組織名	<input type="text"/>
合併日	<input type="text" value="1997"/> 年 <input type="text" value="1"/> 月 <input type="text" value="1"/> 日から <input type="text" value="1997"/> 年 <input type="text" value="12"/> 月 <input type="text" value="31"/> 日まで
検索	

検索結果表示画面				
組織合併情報				
組織名	組織名	新組織名	合併日	見出し

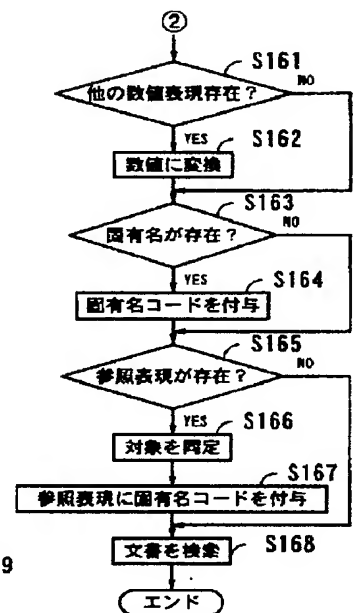
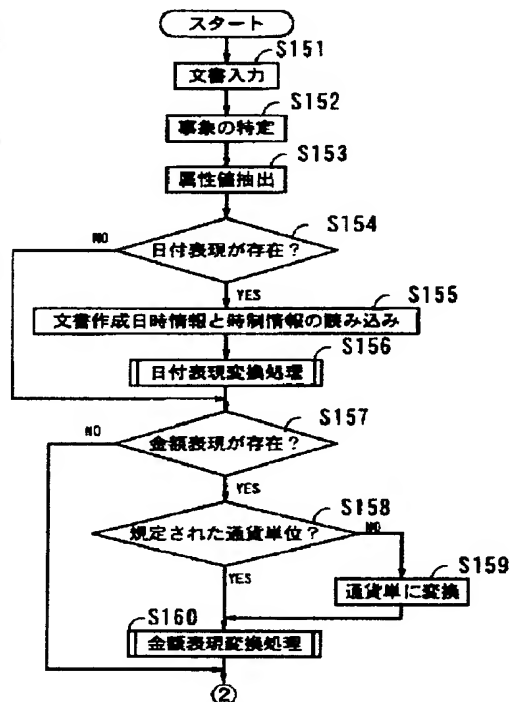


【図26】

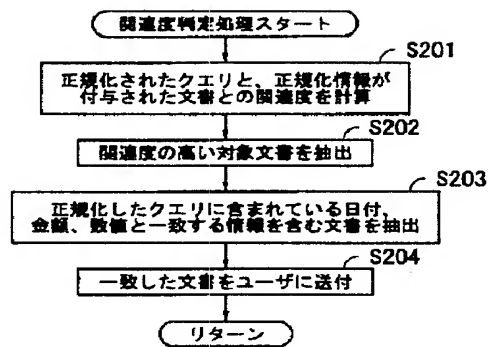
【図27】

【図28】

検索結果表示画面				
組織合併情報				
組織名	組織名	新組織名	合併日	見出し
AAA	BBB	CCC	1997/04/01	AAA, BBB, 2社合併



【図 3 0】



フロントページの続き

(72) 発明者 落谷 亮
神奈川県川崎市中原区上小田中 4 丁目 1 番
1 号 富士通株式会社内

(72) 発明者 西野 文人
神奈川県川崎市中原区上小田中 4 丁目 1 番
1 号 富士通株式会社内